

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-352994

(43)Date of publication of application : 24.12.1999

(51)Int.Cl.

G10L 3/00
G10L 3/00
G10L 3/00
// C12N 15/09

(21)Application number : 10-165030

(71)Applicant : ATR ONSEI HONYAKU TSUSHIN
KENKYUSHO:KK

(22)Date of filing : 12.06.1998

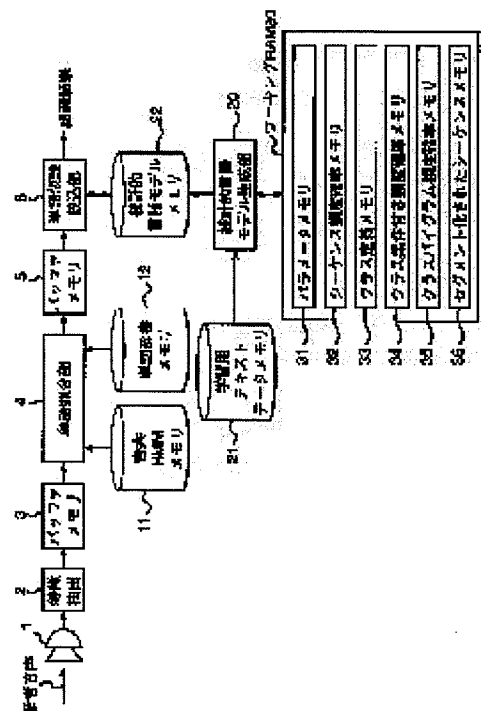
(72)Inventor : SABIN DERIN
KOSAKA YOSHINORI
NAKAJIMA HIDEJI

(54) STATISTICAL SEQUENCE MODEL GENERATOR, STATISTICAL LANGUAGE MODEL GENERATOR, AND SPEECH RECOGNITION SYSTEM

(57)Abstract:

PROBLEM TO BE SOLVED: To impart degree of freedom to an analyzed result, and to process a variable length of sequence as the same class, by re-estimating frequency probability of a bigram between sequences based on front likelihood, frequency probability, and rear likelihood.

SOLUTION: Re-estimation is conducted in a statistical language model generating part 20 to provide the maximum likelihood estimate using EM algorithm based on frequency probability for character lines included in sorted classes and conditional character lines and frequency probability for a bigram between the classes. Frequency probability of a bigram between sequences is re-estimated using an expression for expressing the frequency probability of the bigram between the sequences based on front likelihood for the character lines of processing objects put in a front side of a time series, frequency probability of the character lines when the character line just before the character lines is conditioned, and rear likelihood corresponding to the character lines put in a rear side of the time series, relating to the respective character lines of the processing objects. Bi-multigram statistical sequence models are generated to be output.



(19) 日本国特許庁 (J P)

(12) 特 許 公 報 (B 2)

(11) 特許番号

特許第3004254号
(P3004254)

(45) 発行日 平成12年1月31日 (2000.1.31)

(24) 登録日 平成11年11月19日 (1999.11.19)

(51) Int.Cl. ⁷	識別記号	F I	
G 1 0 L 15/28		G 1 0 L 3/00	5 7 1 D
15/06			5 2 1 C
15/14			5 3 5 Z
// C 1 2 N 15/09		C 1 2 N 15/00	A

請求項の数 9 (全 24 頁)

(21) 出願番号	特願平10-165030	(73) 特許権者	593118597 株式会社エイ・ティ・アール音声翻訳通信研究所 京都府相楽郡精華町大字乾谷小字三平谷5番地
(22) 出願日	平成10年6月12日 (1998.6.12)	(72) 発明者	サビン・デリン 京都府相楽郡精華町大字乾谷小字三平谷5番地 株式会社エイ・ティ・アール音声翻訳通信研究所内
(65) 公開番号	特開平11-352994	(72) 発明者	匂坂 芳典 京都府相楽郡精華町大字乾谷小字三平谷5番地 株式会社エイ・ティ・アール音声翻訳通信研究所内
(43) 公開日	平成11年12月24日 (1999.12.24)	(74) 代理人	100062144 弁理士 青山 葆 (外2名)
審査請求日	平成10年6月12日 (1998.6.12)	審査官	涌井 智則

最終頁に続く

(54) 【発明の名称】 統計的シーケンスモデル生成装置、統計的言語モデル生成装置及び音声認識装置

1

(57) 【特許請求の範囲】

【請求項1】 1個又は複数の単位からなる単位列であるシーケンスを含む入力データに基づいて、可変長の自然数 N_1 個の単位列であるマルチグラムと、可変長の自然数 N_2 個の単位列であるマルチグラムとの間のバイグラムであるパイマルチグラムの統計的シーケンスモデルを生成する統計的シーケンスモデル生成装置であって、

上記入力データに基づいて、予め決められた N_1 、 N_2 の最大値の拘束条件のもとで、すべての単位列の組み合わせの上記バイグラムの頻度確率を計数する初期化手段と、

上記初期化手段によって計数された上記バイグラムの頻度確率に基づいて、各クラスの対をマージしたときの相互情報量の損失が最小となるようにマージして各クラス

2

の頻度確率を更新して予め決められた数の複数のクラスに分類することにより、分類されたクラスに含まれる単位列と、分類されたクラスの条件付きの単位列の頻度確率と、分類されたクラス間のバイグラムの頻度確率を計算して出力する分類手段と、

上記分類処理手段から出力される分類されたクラスに含まれる単位列と、分類されたクラスの条件付きの単位列の頻度確率と、分類されたクラス間のバイグラムの頻度確率とに基づいて、EMアルゴリズムを用いて、最尤推定値を得るように再推定し、ここで、フォワード・バックワードアルゴリズムを用いて、処理対象の各単位列に対して、時系列的に前方にとり得る処理対象の当該単位列に対する前方尤度と、当該単位列の直前の単位列を条件としたときの当該単位列の頻度確率と、時系列的に後方にとり得る当該単位列に対する後方尤度とに基づいて

シーケンス間のバイグラムの頻度確率を示す式を用いて、当該シーケンス間のバイグラムの頻度確率を再推定することにより、再推定結果である上記バイーマルチグラムの統計的シーケンスモデルを生成して出力する再推定手段と、
上記分類手段の処理と上記再推定手段の処理を所定の終了条件を満たすまで繰り返し実行するように制御する制御手段とを備えたことを特徴とする統計的シーケンスモデル生成装置。

【請求項 2】 上記初期化手段はさらに、上記計数されたバイグラムの頻度確率のうち、所定の頻度確率以下のバイグラムの組み合わせのデータを除去することを特徴とする請求項 1 記載の統計的シーケンスモデル生成装置。

【請求項 3】 上記分類手段は、上記初期化手段によって計数された上記バイグラムの頻度確率に基づいて、ブラウンアルゴリズムを用いて、上記複数のクラスに分類することを特徴とする請求項 1 又は 2 記載の統計的シーケンスモデル生成装置。

【請求項 4】 上記式は、上記入力データにおいて、当該単位列である第 2 の単位列が第 1 の単位列に続くときの単位列のシーケンス間のバイグラムの頻度確率を、上記入力データにおける処理対象の各単位列に対して計算するための式であり、
上記シーケンス間のバイグラムの頻度確率は、第 1 と第 2 の単位列を含むすべてのセグメント化での尤度の和を、第 1 の単位列を含むすべてのセグメント化での尤度の和で除算することによって得られたことを特徴とする請求項 1 乃至 3 のうちの 1 つに記載の統計的シーケンスモデル生成装置。

【請求項 5】 上記式は、上記入力データにおいて各単位列が発生する平均回数を示す分母と、上記入力データにおいて第 2 の単位列が第 1 の単位列に続くときの各単位列に対する平均回数を示す分子とを有し、
上記分子は、処理対象の各単位列に対する、上記前方尤度と、当該単位列の直前の単位列を条件としたときの当該単位列の頻度確率と、上記後方尤度の積の和であり、
上記分母は、処理対象の各単位列に対する、上記前方尤度と、当該単位列の直前の単位列を条件としたときのすべての単位列の頻度確率と、上記後方尤度の積の和であることを特徴とする請求項 4 記載の統計的シーケンスモデル生成装置。

【請求項 6】 上記終了条件は、上記分類手段の処理と、上記再推定手段の処理との反復回数が予め決められた回数に達したときであることを特徴とする請求項 1 乃至 5 のうちの 1 つに記載の統計的シーケンスモデル生成装置。

【請求項 7】 請求項 1 乃至 6 のうちの 1 つに記載の統計的シーケンスモデル生成装置において、
上記単位は自然言語の文字であり、上記シーケンスは単

語であり、上記分類手段は、文字列を複数の単語の列に分類し、上記統計的シーケンスモデルは、統計的言語モデルであることを特徴とする統計的言語モデル生成装置。

【請求項 8】 請求項 1 乃至 6 のうちの 1 つに記載の統計的シーケンスモデル生成装置において、上記単位は自然言語の単語であり、上記シーケンスはフレーズであり、上記分類手段は、単語列を複数のフレーズの列に分類し、上記統計的シーケンスモデルは、統計的言語モデルであることを特徴とする統計的言語モデル生成装置。

【請求項 9】 入力される発声音声文の音声信号に基づいて、所定の統計的言語モデルを用いて音声認識する音声認識手段を備えた音声認識装置において、上記音声認識手段は、請求項 7 又は 8 記載の統計的言語モデル生成装置によって生成された統計的言語モデルを参照して音声認識することを特徴とする音声認識装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、学習用シーケンスデータに基づいて統計的シーケンスモデルを生成する統計的シーケンスモデル生成装置、学習用テキストデータに基づいて統計的言語モデルを生成する統計的言語モデル生成装置、及び上記統計的言語モデルを用いて、入力される発声音声文の音声信号を音声認識する音声認識装置に関する。

【0002】

【従来の技術】近年、連続音声認識装置において、その性能を高めるために言語モデルを用いる方法が研究されている。これは、シーケンスモデルである言語モデルを用いて、次単語を予測し探索空間を削減することにより、認識率の向上及び計算時間の削減の効果を狙ったものである。ここで、シーケンスとは、具体的には、文字のシーケンスでは単語であり、単語のシーケンスではフレーズ（又は句）である。最近盛んに用いられている言語モデルとして N-gram（N-グラム；ここで、N は 2 以上の自然数である。）がある。これは、大規模なテキストデータを学習し、直前の N-1 個の単語から次の単語への遷移確率を統計的に与えるものである。複数 L 個の単語列 $w_1^L = w_1, w_2, \dots, w_L$ の生成確率 $P(w_1^L)$ は次式で表される。

【0003】

【数 1】

$$P(w_1^L) = \prod_{t=1}^L P(w_t | w_{t-1-n}^{t-1})$$

【0004】ここで、 w_t は単語列 w_1^L のうち t 番目の 1 つの単語を表し、 w_i^j は i 番目から j 番目の単語列を表す。上記数 1 において、確率 $P(w_t | w_{t-1-n}^{t-1})$ は、N 個の単語からなる単語列 w_{t-1-n}^{t-1}

が発声された後に単語 w_t が発声される確率であり、以下同様に、確率 $P(A|B)$ は単語又は単語列 B が発声された後に単語 A が発声される確率を意味する。また、数1における「 Π 」は $t=1$ から L までの確率 $P(w_t | w_{1:t-1})$ の積を意味する。

【0005】ところで、近年、上記統計的言語モデルの N -gram を用いて連続音声認識の性能を向上させる手法が盛んに提案されており、そのいくつかのモデルでは、可変長の単語列にわたる単語の依存性を利用する方法を用いている。これらのモデルは、共通して従来の N -gram モデルにみられる固定長の依存性の仮定を緩和するために用いられており、種々のより広い仮定をカバーしている。

【0006】フレーズを純粹に統計的方法（すなわち、統計的文脈自由文法 (Stochastic Context Free Grammars) にあるような文法的規則を用いない方法）で導くためには、種々の基準を使用する必要がある、例えば、以下の基準が提案されてきた。

(a) 従来技術文献1「K. Ries et al., "Class phrase models for language modeling", Proceedings of ICSLP 96, 1996」において開示されたリーブ・ワン・アウト (leave-one-out) 尤度、及び

(b) 従来技術文献2「H. Masataki et al., Variable-order n -gram generation by word-class splitting and consecutive word grouping. Proceedings of ICASSP 96, 1996」において開示されたエントロピー。

【0007】

【発明が解決しようとする課題】これらの方法において、尤度の基準を統計的枠組みの中で用いることで、EM (Expectation Maximum; すなわち、期待値の最大化) アルゴリズムを用いた最適化の方法を用いることができるが、過学習となる傾向がある。また、最適化処理においては、例えば、従来技術文献3「S. Matsunaga et al., "Variable-length language modeling integrating global constraints", Proceedings of EUROSpeech 97, 1997」において発見的手法を用いられているが、統計的言語モデルの収束と最適化は理論的に保証されていない。

【0008】ここで、さらに、例えば、従来技術文献1において提案された尤度の基準を用いたときの問題点について述べると以下の通りである。

<問題点1>単語のシーケンスの頻度確率が貪欲なアルゴリズム (greedy algorithm) によって得られるために、最適状態に向かう単調な収束が保証されない。

<問題点2>この方法は確定的なものである。つまり、仮にシーケンス $[bcd]$ がシーケンスの目録 (inventory) に在れば、入力文字列に $"bcd"$ が発生しても、これが $[bc] + [d]$ 、 $[b] + [cd]$ 、 $[b] + [c] + [d]$ 等のサブシーケンスに分割されることはない。言い換えれば、シーケンスへの解析において自由

度が無い。

<問題点3>シーケンスのクラスの定義が先行する単語のクラス分類を基礎としている。すなわち、まず、単語が分類され、次に、単語のクラスのラベルの各シーケンスは、シーケンスのクラスを定義するために使用される。従って、同一クラスに長さの違うシーケンスを入れることはできない。例えば、 $"thank\ you\ for"$ と $"thank\ you\ very\ much\ for"$ は同じクラスに入らない。

【0009】これを解決するために、本発明者は、従来技術文献4「S. Deligne et al., "Introducing statistical dependencies and structural constraints in variable-length sequence models", In Grammatical Inference: Learning Syntax from Sentences, Lecture Notes in Artificial Intelligence 1147, pp. 156-167, Springer, 1996」において、可変長のシーケンスであるマルチグラムを用いる統計的言語モデルについて、当該従来技術文献4の(16)式を用いて、それらのパラメータを計算できる可能性だけを示しているが、当該(16)式は、実際にディジタル計算機を用いて計算することができる形式とはなっておらず、実用化することができないという問題点があった。ここで、マルチグラムとは、他のシーケンスとの依存性を特定しない可変長のシーケンスである。

【0010】本発明の目的は以上の問題点を解決し、従来例に比較して、最適な状態に向かう単調な収束を保証することができ、解析結果に自由度があり、可変長のシーケンスを同一のクラスで取り扱うことができ、ディジタル計算機を用いて実用的に高速処理して統計的モデルを生成することができる統計的シーケンスモデル生成装置、統計的言語モデル生成装置及び音声認識装置を提供することにある。

【0011】

【課題を解決するための手段】本発明に係る統計的シーケンスモデル生成装置は、1個又は複数の単位からなる単位列であるシーケンスを含む入力データに基づいて、可変長の自然数 N_1 個の単位列であるマルチグラムと、可変長の自然数 N_2 個の単位列であるマルチグラムとの間のバイグラムであるバイーマルチグラムの統計的シーケンスモデルを生成する統計的シーケンスモデル生成装置であって、上記入力データに基づいて、予め決められた N_1 、 N_2 の最大値の拘束条件のもとで、すべての単位列の組み合わせの上記バイグラムの頻度確率を計数する初期化手段と、上記初期化手段によって計数された上記バイグラムの頻度確率に基づいて、各クラスの対をマージしたときの相互情報量の損失が最小となるようにマージして各クラスの頻度確率を更新して予め決められた数の複数のクラスに分類することにより、分類されたクラスに含まれる単位列と、分類されたクラスの条件付きの単位列の頻度確率と、分類されたクラス間のバイグラムの頻度確率を計算して出力する分類手段と、上記分類処

理手段から出力される分類されたクラスに含まれる単位列と、分類されたクラスの条件付きの単位列の頻度確率と、分類されたクラス間のバイグラムの頻度確率とに基づいて、EMアルゴリズムを用いて、最尤推定値を得るように再推定し、ここで、フォワード・バックワードアルゴリズムを用いて、処理対象の各単位列に対して、時系列的に前方にとり得る処理対象の当該単位列に対する前方尤度と、当該単位列の直前の単位列を条件としたときの当該単位列の頻度確率と、時系列的に後方にとり得る当該単位列に対する後方尤度とに基づいてシーケンス間のバイグラムの頻度確率を示す式を用いて、当該シーケンス間のバイグラムの頻度確率を再推定することにより、再推定結果である上記バイマルチグラムの統計的シーケンスモデルを生成して出力する再推定手段と、上記分類手段の処理と上記再推定手段の処理を所定の終了条件を満たすまで繰り返し実行するように制御する制御手段とを備えたことを特徴とする。

【0012】また、上記統計的シーケンスモデル生成装置において、上記初期化手段はさらに、上記計数されたバイグラムの頻度確率のうち、所定の頻度確率以下のバイグラムの組み合わせのデータを除去することを特徴とする。

【0013】さらに、上記統計的シーケンスモデル生成装置において、上記分類手段は、上記初期化手段によって計数された上記バイグラムの頻度確率に基づいて、ブラウンアルゴリズムを用いて、上記複数のクラスに分類することを特徴とする。

【0014】また、上記統計的シーケンスモデル生成装置において、上記式は、上記入力データにおいて、当該単位列である第2の単位列が第1の単位列に続くときの単位列のシーケンス間のバイグラムの頻度確率を、上記入力データにおける処理対象の各単位列に対して計算するための式であり、上記シーケンス間のバイグラムの頻度確率は、第1と第2の単位列を含むすべてのセグメント化での尤度の和を、第1の単位列を含むすべてのセグメント化での尤度の和で除算することによって得られる。また、ここで、上記式は、上記入力データにおいて各単位列が発生する平均回数を示す分母と、上記入力データにおいて第2の単位列が第1の単位列に続くときの各単位列に対する平均回数を示す分子とを有し、上記分子は、処理対象の各単位列に対する、上記前方尤度と、当該単位列の直前の単位列を条件としたときの当該単位列の頻度確率と、上記後方尤度の積の和であり、上記分母は、処理対象の各単位列に対する、上記前方尤度と、当該単位列の直前の単位列を条件としたときのすべての単位列の頻度確率と、上記後方尤度の積の和である。

【0015】さらに、上記統計的シーケンスモデル生成装置において、上記終了条件は、上記分類手段の処理と、上記再推定手段の処理との反復回数が予め決められた回数に達したときであることを特徴とする。

【0016】また、本発明に係る統計的言語モデル生成装置は、上記統計的シーケンスモデル生成装置において、上記単位は自然言語の文字であり、上記シーケンスは単語であり、上記分類手段は、文字列を複数の単語の列に分類し、上記統計的シーケンスモデルは、統計的言語モデルであることを特徴とする。

【0017】さらに、本発明に係る統計的言語モデル生成装置は、上記統計的シーケンスモデル生成装置において、上記単位は自然言語の単語であり、上記シーケンスはフレーズであり、上記分類手段は、単語列を複数のフレーズの列に分類し、上記統計的シーケンスモデルは、統計的言語モデルであることを特徴とする。

【0018】またさらに、本発明に係る音声認識装置は、入力される発声音声文の音声信号に基づいて、所定の統計的言語モデルを用いて音声認識する音声認識手段を備えた音声認識装置において、上記音声認識手段は、上記統計的言語モデル生成装置によって生成された統計的言語モデルを参照して音声認識することを特徴とする。

【0019】

【発明の実施の形態】以下、図面を参照して本発明に係る実施形態について説明する。以下の実施形態においては、単位は文字であり、文字のシーケンスである文字列を単語列に分類する一例、並びに、単位は単語であり、単語のシーケンスである単語列をフレーズ（句）に分類する一例について説明しているが、本発明はこれに限らず、単位はDNAであり、DNAのシーケンスであるDNA列を所定のDNA配列に分類するように構成してもよい。また、単位は塩基であり、塩基のシーケンスである塩基列を所定のコドンに分類するように構成してもよい。

【0020】図1は、本発明に係る一実施形態である連続音声認識装置のブロック図である。本実施形態の連続音声認識装置は、学習用テキストデータメモリ21に記憶された文字列であるテキストデータに基づいて、ワーキングRAM30を用いて、可変長のバイマルチグラムの言語モデルを生成する統計的言語モデル生成部20を備え、ここで、統計的言語モデル生成部20の処理は、図3に示すように、大きく分けると、ブラウンアルゴリズムを用いた分類処理（ステップS3）と、バイマルチグラムを用いた再推定処理（ステップS4）とを含むことを特徴としている。

【0021】すなわち、本実施形態の統計的言語モデル生成装置は、1個又は複数の文字からなる文字列のシーケンスを含む入力データに基づいて、可変長の自然数 N_1 個の文字列と可変長の自然数 N_2 個の文字列との間のバイグラムであるバイマルチグラムの統計的言語モデルを生成する統計的言語モデル生成装置であり、ここで、図3に示すように、（a）上記入力データに基づいて、予め決められた N_1 、 N_2 の最大値の拘束条件のもとで、

すべての文字列の組み合わせの上記バイグラムの頻度確率を計数する初期化処理（ステップS2）と、（b）上記初期化処理によって計数された上記バイグラムの頻度確率に基づいて、各クラスの対をマージしたときの相互情報量の損失が最小となるようにマージして各クラスの頻度確率を更新して予め決められた数の複数のクラスに分類することにより、分類されたクラスに含まれる文字列と、分類されたクラスの条件付きの文字列の頻度確率と、分類されたクラス間のバイグラムの頻度確率を計算して出力する分類処理（ステップS3）と、（c）上記分類処理によって得られた分類されたクラスに含まれる文字列と、分類されたクラスの条件付きの文字列の頻度確率と、分類されたクラス間のバイグラムの頻度確率とに基づいて、EMアルゴリズムを用いて、最尤推定値を得るように再推定し、ここで、フォワード・バックワードアルゴリズムを用いて、処理対象の各文字列に対して、時系列的に前方にとり得る処理対象の当該文字列に対する前方尤度と、当該文字列の直前の文字列を条件としたときの当該文字列の頻度確率と、時系列的に後方にとり得る当該文字列に対する後方尤度とに基づいてシーケンス間のバイグラムの頻度確率を示す式（数22-数24）を用いて、当該シーケンス間のバイグラムの頻度確率を再推定することにより、再推定結果である上記バイマルチグラムの統計的シーケンスモデルを生成して出力する再推定処理（ステップS4）と、（d）上記分類処理と上記再推定処理を所定の終了条件を満たすまで繰り返し実行するように制御する処理（ステップS5）を含むことを特徴とする。

【0022】本実施形態では、単語のN-gramに基づく手法に対向する、フレーズに基づく方法に焦点を当てる。ここで、複数の文はフレーズに構成され、頻度確率は、単語に代わってフレーズに割り当てられる。モデルがN-gramに基づくか、フレーズに基づくかに関わらず、それらは確定的モデルあるいは統計的モデルのいずれかに該当する。フレーズに基づく枠組みでは、非確定性はその文の解析結果の曖昧さを通じてフレーズに導入される。すなわち、これは実際においては、フレーズ“abc”がフレーズとして登録されているにもかかわらず、文字列の解析結果が例えば[ab][c]となる確率が皆無でないことを意味する。これとは対照的に、確定的手法ではa、b、cすべての同時出現はシステマティックにフレーズ[abc]の出現と解釈される。

【0023】また、本実施形態では、統計的言語モデルの処理は、バイマルチグラムを用いて実行され、当該バイマルチグラムの言語モデルは、フレーズに基づく統計的モデルであり、そのパラメータは尤度基準に従って推定される。

【0024】まず、マルチグラムの理論的な定式化につ*
尤度

*いて説明する。マルチグラムの枠組みでは、T個の単語からなる文

【数2】 $W = w^{(1)} w^{(2)} \cdots w^{(T)}$

は、それぞれ最大長n個の単語からなる各々のフレーズが連鎖（シーケンス）したものと仮定される。ここで、SはT個のフレーズへのセグメント化を示し、 $s^{(t)}$ はセグメント化Sにおける時刻インデックス（最初の語からのシリアル番号を示す。）（t）のフレーズとした場合、WのSでのセグメント化の結果は、次式で表すことができる。

【数3】 $(W, S) = s^{(1)} \cdots s^{(Ta)}$

【0025】ここで、セグメント化された複数のフレーズからなる辞書は、語彙から1, 2...からnにいたるまでの単語を組み合わせ形成されるものであり、ここでは、次式のように表す。

【数4】 $D_s = \{s_j\}_j$

そして、文の尤度は、各セグメント化に対する尤度の和として、次式のように計算される。

【0026】

【数5】

$$L(W) = \sum_{S \in \{S\}} L(W, S)$$

【0027】モデルの決定指向的手法により、文Wは、最も尤らしいセグメント化に従って解析され、次の近似式が得られる。

【0028】

【数6】

$$L^*(W) = \max_{S \in \{S\}} L(W, S)$$

【0029】ここで、フレーズ間のn-gramの相関を仮定し、特定のセグメント化Sの結果の尤度の値を次式のように計算する。

【0030】

【数7】

$$L(W, S) = \prod_t p(s^{(t)} | s^{(t-n+1)} \cdots s^{(t-1)})$$

【0031】ここで、以下、符号nは複数のフレーズ間の依存度を表し、従来のn-gramの表記法のnとして使用する。また、符号 n_{\max} は、フレーズの最大長を表す。従って、ここで、尤度の計算例を次式に示す。この例では、バイマルチグラムモデル（ $n_{\max} = 3$, $n = 2$ ）の“abcd”の尤度を示す。記号#は空のシーケンスを表す。

【0032】

【数8】

11

12

$$\begin{aligned}
&= p([a] | \#) p([b] | [a]) p([c] | [b]) p([d] | [c]) \\
&+ p([a] | \#) p([b] | [a]) p([cd] | [b]) \\
&+ p([a] | \#) p([bc] | [a]) p([d] | [bc]) \\
&+ p([a] | \#) p([bcd] | [a]) \\
&+ p([ab] | \#) p([c] | [ab]) p([d] | [c]) \\
&+ p([ab] | \#) p([cd] | [ab]) \\
&+ p([abc] | \#) p([d] | [abc])
\end{aligned}$$

【0033】上記数8から明らかなように、当該尤度は、シーケンス“abcd”をセグメント化するときのすべての組み合わせについての頻度確率の和を表している。

【0034】次いで、言語モデルのパラメータの推定について説明する。マルチグラムのn-gramモデルは、パラメータ θ のセットによって完全に定義され、次式のパラメータ θ は、辞書Dsを用いて、

【数9】

$$\theta = \{ p(s_{in} | s_{i1} \dots s_{in-1}) | s_{i1} \dots s_{in} \in Ds \}$$

n個のフレーズのあらゆる組み合わせに関係するn-gramの条件付き確率によって構成される。パラメータ*

$$Q(k, k+1) = \sum_{S \in \{S\}} L^{(k)}(S | W) \log \{ L^{(k+1)}(W, S) \}$$

【0036】公知のEMアルゴリズムにおいて示されるように、

$$\text{【数11】 } Q(k, k+1) \geq Q(k, k)$$

であれば、

$$\text{【数12】 } L^{(k+1)}(W) \geq L^{(k)}(W)$$

である。従って、反復回数パラメータ(k+1)における次式の再推定式

$$\text{【数13】 } p^{(k+1)}(s_{in} | s_{i1} \dots s_{in-1})$$

は、次式の拘束条件

【数14】

$$p^{(k+1)}(s_{in} | s_{i1} \dots s_{in-1})$$

$$= p_i / p_i$$

ここで、

$$p_i = \{ \sum_{S \in \{S\}} c(s_{i1} \dots s_{in-1} s_{in}, S) \times L^{(k)}(S | W) \}$$

$$p_b = \{ \sum_{S \in \{S\}} c(s_{i1} \dots s_{in-1}, S) \times L^{(k)}(S | W) \}$$

【0038】ここで、 $c(s_{i1} \dots s_{in}, S)$ は、セグメント化Sにおける複数のフレーズ $s_{i1} \dots s_{in}$ の組み合わせの出現数を示す。数15の再推定式は、バイーマルチグラム(n=2)について詳細後述されるように、フォワード・バックワードアルゴリズム(forward backward algorithm)(以下、FB法ともいう。)を用いて実行される。決定指向の方法では、再推定式は、次式のように簡略化される。

【0039】

$$\text{【数16】 } p^{(k+1)}(s_{in} \dots s_{in-1}) = \{ c(s_{i1} \dots s_{in-1}, S^{*(k)}) \} / \{ c(s_{i1} \dots s_{in-1}, S^{*(k)}) \}$$

* θ のセットの推定値は、例えば、不完全なデータから得られる想定しうる最大の尤度値、すなわち最尤推定値(Maximum Likelihood Estimation)として得られ、ここで、未知のデータは基礎をなすセグメント化Sである。従って、パラメータ θ の反復的な最尤推定値は、公知のEMアルゴリズム(Expectation Maximization Algorithm)によって計算することができる。ここで、 $Q(k, k+1)$ を、反復回数パラメータk及びk+1の尤度を用いて計算される、次式の補助関数とする。

【0035】

【数10】

$$\sum_{s_{in} \in D_s} p(s_{in} | s_{i1} \dots s_{in-1}) = 1$$

のもとで、モデルパラメータ $\theta^{(k+1)}$ について補助関数 $Q(k, k+1)$ を最大化することにより、次式のように直接的に導くことができる。なお、本明細書において、下付きの下付きの表記及び上付きの下付きの表記はできないので、下層の下付きの表記を省略している。

【0037】

【数15】

$$\{ c(s_{i1} \dots s_{in-1}, S^{*(k)}) \} / \{ c(s_{i1} \dots s_{in-1}, S^{*(k)}) \}$$

【0040】ここで、 $S^{*(k)}$ は、 $L^{(k)}(S | W)$ を最大化する文の解析結果であり、ビタビ(Viterbi)アルゴリズムによって導かれる。各反復は、尤度 $L^{(k)}(W)$ を増大させる意味において言語モデルを改善し、最終的には臨界点(おそらくは、局所最大値)へ収束する。モデルパラメータ θ のセットは、学習用コーパス、すなわち学習用テキストデータにおいて観察されるすべてのフレーズの組み合わせの相対的頻度を用いて初期化され

る。

【0041】次いで、可変長フレーズのクラスタリング（分類処理）について説明する。従来技術文献1によれば、近年、クラスフレーズに基づくモデルが注目されているが、通常、それは従来の単語クラスタリングを仮定している。典型的には、各単語はまず、単語が属するクラスのラベル C_k を割り当てられ、単語ークラスラベルの可変長フレーズ $[C_{k1}, C_{k2} \dots C_{kn}]$ が導かれる。各可変長フレーズによって、“ $< [C_{k1}, C_{k2} \dots C_{kn}]$ ”として示されるフレーズが属するクラスのラベルが定義される。しかしながら、この手法では、同じ長さのフレーズのものにしか同じフレーズークラスラベルを割り当てることができない。例えば、“thank you for”と“thank you very much for”というフレーズを同じクラスラベルに割り当てることができない。本実施形態では、このような限界に対する解決法として、単語に代わり直接フレーズをクラスタリングする方法を提案する。この目的を達成するためには、2個のフレーズ間のバイグラム（ $n_{\max} = 2$ ）を仮定し、上述したバイマルチグラムモデルの学習手法に変更を加え、各反復が次の2つの段階より構成されるようにする。

【0042】(I) ステップSS1：クラス割り当て（図3のステップS3に対応する。）

$L(W, S)$

$$= \prod_t p(C_{k(t-1)} | C_{k(t-2)}) p(s_{(t)} | C_{k(t-1)})$$

【0045】これは、上述したように、頻度確率 $p^{(k)}(s_j | s_i)$ に対する処理と同様に、頻度確率 $p^{(k)}(C_{k(s_j)} | C_{k(s_i)}) \times p^{(k)}(s_j | C_{k(s_j)})$ に基づいて頻度確率 $p^{(k+1)}(s_j | s_i)$ を再推定することに等しい。

【0046】要約すれば、上記ステップSS1によって、現在のフレーズ分布に関し、相互情報量の基準に基づくクラス割り当てが最適化されるよう保証され、上記ステップSS2によって、現在のクラスの頻度確率を用いて、上記数19に従って、計算された尤度がフレーズの頻度確率により最適化されるよう保証される。学習データは、従って、完全に統合化された方法により連合的（paradigmatic）かつ統合的（syntagmatic）（それぞれ言語学の用語である。）レベルの双方において反復的に構成される。すなわち、クラス割り当てにより表現されるフレーズ間の連合的關係はフレーズの頻度確率の再推定に影響を与え、フレーズの頻度確率は後続するクラス割り当てを決定する。

【0047】本実施形態では、上述のように、バイマルチグラムのパラメータの推定のために、フォワード・バックワードアルゴリズム（FB法）を用いる。これについて、以下に、詳述する。

【0048】上記数15は、フォワード・バックワード

* 【数17】 $\{p^{(k)}(s_j | s_i)\} \rightarrow \{p^{(k)}(C_{k(s_j)} | C_{k(s_i)}), p^{(k)}(s_j | C_{k(s_i)})\}$

(II) ステップSS2：マルチグラムの再推定（図3のステップS4に対応する。）

【数18】 $\{p^{(k)}(C_{k(s_j)} | C_{k(s_i)}), p^{(k)}(s_j | C_{k(s_i)})\} \rightarrow \{p^{(k+1)}(s_j | s_i)\}$

【0043】上記ステップSS1では、フレーズバイグラムの頻度確率を入力とし、クラスバイグラムの頻度確率を出力する。クラス割り当ては、例えば、従来技術文献5「P. F. Brown et al., "Class-based n-gram models of natural language", Computational Linguistics, Vol. 18, No. 4, pp. 467-479, 1992」によれば、隣り合うフレーズ間の相関情報を最大化することによって行われる。ここで、クラスタリングの候補は単語ではなくフレーズとする。上述のように、 $\{p^{(k)}(s_j | s_i)\}$ は、学習用テキストデータにおけるフレーズの同時出現の相対的頻度を用いて初期化される。上記ステップSS2では、マルチグラムの再推定式（数15）又はその近似式（数16）を用いてフレーズの頻度確率を再推定する。ここで、唯一の違いは、解析結果の尤度は以下の式により計算される。

【0044】

* 【数19】

アルゴリズムを用いて、 n_{\max} をシーケンスの最大長とし、 T をコーパス（学習用テキストデータ）の語数として、複雑さの度合いであるコンプレキシティ $O(n_{\max}^2 T)$ で計算することができる。ここで、コンプレキシティ $O(n_{\max}^2 T)$ は計算コストのオーダーに対応する。すなわち、当該数15の計算コストは、シーケンスの最大長 n_{\max} の2乗に比例し、コーパスの語数に比例する。本実施形態においては、基本的には、セグメント化 $\{S\}$ のセットではなく、単語のタイムインデックス (t) にわたって加算を行い、数15の分子及び分母を計算する。ここで、当該計算は、次式の前方向の変数 $\alpha(t, l_i)$ 及び後ろ方向の変数 $\beta(t, l_j)$ の定義に依存する。

【0049】

【数20】

$$\alpha(t, l_i) = L(W_{(t)}^{(t-1)} | [W_{(t-1)}^{(t)}])$$

【数21】

$$\beta(t, l_j) = L(W_{(t+1)}^{(t)} | [W_{(t-j+1)}^{(t)}])$$

【0050】前方向の変数 $\alpha(t, l_i)$ は、最初の t 個の単語の尤度を表し、ここで、最後の l_i 個の単語は、1つのシーケンスを形成するように制限される。また、後ろ方向の変数 $\beta(t, l_j)$ は、最後の $(T-t)$ 個

の語の条件付き尤度を示し、最後の $(T-t)$ 個の単語は、シーケンス $[w_{(t-1)+1} \dots w_{(t)}]$ に後続する。ここで、例えば、 $w_{(t)}$ は、時刻インデックス (1) から $(t-1)$ までの単語からなる単語列を表す。そして、解析結果の尤度は、数7によって計算されると仮定すると、数15は次式のように書き換えられ

$$p_t = \sum_{\mathbf{T}} \alpha(t, l_1) p^{(t)}(s_j | s_i) \beta(t+1, l_1) \delta_i(t-1+1) \delta_j(t+1) \\ t=1$$

【数24】

$$p_t = \sum_t \alpha(t, l_1) \beta(t, l_1) \delta_i(t-1+1)$$

【0052】ここで、 l_1 及び l_j はそれぞれシーケンス s_i 及び s_j の長さを示す。クロネッカー関数 $\delta_k(t)$ は、時刻インデックス t で開始する単語のシーケンスが s_k であるときは1となる一方、そうでない場合は0となる関数である。また、変数 α 及び β は以下の反復式

$$\alpha(t, l_1) = \sum_{l=1}^{n_{\max}} \alpha(t-1, l) p([w_{(t-1)+1}^{(t)}] | [w_{(t-1)+1}^{(t-1)}])$$

ここで、

$$\alpha(0, 1) = 1, \alpha(0, 2) = \dots = \alpha(0, n_{\max}) = 0$$

$$\beta(t, l_1)$$

$$= \sum_{l=1}^{n_{\max}} p([w_{(t+1)}^{(t+1)}] | [w_{(t-1)+1}^{(t)}]) \beta(t+1, l_1)$$

ここで、

$$\beta(T+1, 1) = 1, \beta(T+1, 2) = \dots = \beta(T+1, n_{\max}) = 0$$

である。

【0055】解析結果の尤度がクラスの仮定を用いて計算される場合、すなわち、数19に従って計算される場合は、再推定式(数22-数24)の項 $p^{(t)}(s_j | s_i)$ はそのクラスの等価物、すなわち $p^{(t)}(C_{k(s_j)} | C_{k(s_i)})$ に置き換えられる。 α の反復式において、項 $p([w_{(t-1)+1}^{(t)}] | [w_{(t-1)+1}^{(t-1)}])$ は、シーケンス $[w_{(t-1)+1}^{(t)}]$ のクラスの条件付き確率を乗算した

対応するクラスのバイグラム確率に置き換えられる。同様の変形を反復式における変数 β についても行う。

【0056】次いで、本実施形態におけるフォワード・バックワードアルゴリズムを用いた再推定処理について、一例を参照して、以下に詳述する。前方向及び後

ろる。

$$[0051]$$

$$[\text{数}22] p^{(k+1)}(s_j | s_i) = p_c / p_d$$

ここで、

$$[\text{数}23]$$

※(又は帰納式)によって計算できる。ここで、時刻インデックス $t=0$ 及び $t=T+1$ においてそれぞれ開始及び終了シンボルを仮定する。

$$[\text{数}25] 1 \leq t \leq T+1 \text{ に対して:}$$

★である。

$$[\text{数}27] 0 \leq t \leq T \text{ に対して:}$$

方向(以下、前後方向という。)の再推定処理は、数22の分子の加算、及び分母の加算が、可能な解析結果集合 $\{S\}$ に代わって、学習データにおける単位の時刻インデックス t について計算されるように、数15における複数の項を配列し直して行う。この方法は、前方向の変数 α 及び後ろ方向の変数 β の定義に依存している。

(a) 下記のパラグラフ<<A1>>では、クラスのなにことを仮定している。

(b) 下記のパラグラフ<<A1.1>>では、変数 α 及び β を定義し、例を提供する。

(c) 下記のパラグラフ<<A1.2>>では、変数 α 及び β を使用した頻度確率に関する前後方向の再推定について例示する。

(d) 下記のパラグラフ<<A1.3>>では、反復(又は帰納)による変数 α と β の計算方法に関して例示する。

(e) 下記のパラグラフ<<A2>>では、クラスが存

在する場合のパラグラフ<<A1. 2>>及び<<A

1. 3>>の修正方法を示す。

(f) 下記の例はすべて、次の表に示すデータに基づい*

* ている。

【0057】

【表1】

入学習データ(下記) :

onesixone e i g h t s i x t h r e e t w o
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25

単位の時刻インデックス(上記) :

(注) 学習データの1つの文字は、1つの時刻インデックスに対応している。

【0058】<<A1. 1>>前方向の変数 α 及び後ろ方向の変数 β の定義

変数 $\alpha(t, 1)$ は、長さ1のシーケンスで終了する、時刻インデックス(t)までのデータの尤度である。例えば、変数 $\alpha(9, 3)$ は、シーケンス「one six one」の尤度である。また、変数 $\beta(t, 1)$ は、長さ1のシーケンスが時刻インデックス(t)で終了するということが知られているときに、時刻インデックス(t+1)で開始されるデータの条件つき尤度である。例えば、変数 $\beta(9, 3)$ は、先行するシーケンスが「one 20 e」であるときの、シーケンス「eight six three two」の尤度である。反復又は帰納による変数 α 及び β の計算方法に関する例を、下記のパラグラフ<<A1. 3>>に示す。

【0059】<<A1. 2>>変数 α 及び β に基づく確率の再推定

例として、上記の学習データ例に関する、変数 α 及び β を使用した頻度確率 $p(o_n_e | s_i_x)$ の再推定式を示す。頻度確率 $p(o_n_e | s_i_x)$ の一般的な再推定式(数15))は次のような意味を持つ。

(a) 分子は、学習データにおいてシーケンス「one」がシーケンス「six」に続く平均回数である。

(b) 分母は、学習データにおいてシーケンス「six」が発生する平均回数である。

(c) ここで、平均回数の値は、学習データのシーケンスにおけるすべての可能な解析結果について求める。

【0060】フォワード・バックワードアルゴリズムを用いた再推定式(数22-24)の分子(数23)及び分母(数24)はそれぞれ、数15の分子及び分母に等しいが、これらは解析結果集合にわたる加算ではなく、時刻インデックスにわたる加算によって計算したものである。再推定式(数15)の分子では、「six」と「one」の2つのシーケンスが連続して発生する毎に、各可能な解析結果の尤度が加算される。一方、フォワード・バックワードアルゴリズムを用いた再推定式(数22-数24)においては、「six」と「one」の2つのシーケンスが連続して発生し、また、シーケンス「one」が時刻インデックス(t+1)で開始するようすすべての解析結果の尤度値をまずグループ化して、加算する。時刻インデックスtまで加算した時点で加算計算は 50

完了する。

【0061】上記の例では、「six」と「one」の2つのシーケンスが連続して発生し、しかもシーケンス「one」が時刻インデックス(7)でのみ開始している。ここで、「six」と「one」の2つのシーケンスが連続して発生し、また、時刻インデックス(7)でシーケンス「one」が開始するようすすべての解析結果の尤度値の和は、シーケンス「one six one e i g h t s i x t h r e e t w o」の尤度であり、これは、次式に等しい。

【数29】

(k)

$$\alpha(6, 3) \times p(o_n_e | s_i_x) \times \beta(9, 3)$$

【0062】ここで、第2項の $p(o_n_e | s_i_x)$ は、反復回数パラメータ(k)における頻度確率である。また、前方向の変数 α の定義により、変数 $\alpha(6, 3)$ はシーケンス「one six」の尤度であり、さらに、後ろ方向の変数 β の定義により、変数 $\beta(9, 3)$ は、シーケンス「one」が得られたときの、シーケンス「eight six three two」の尤度である。

【0063】数15の分母では、可能な各解析結果の尤度を、シーケンス「six」がこの解析において発生するのと同じ回数で加算する。等価である、フォワード・バックワードアルゴリズムを用いた前後方向の定式化では、シーケンス「six」が発生し、時刻インデックス(t)で終了するすべての全解析結果の尤度値をまずグループ化した後に加算し、時刻インデックスtを越えた時点で加算を終了する。

【0064】上述の例では、シーケンス「six」は、時刻インデックス(6)と時刻インデックス(17)で終了するように発生している。シーケンス「six」が時刻インデックス(6)で終了するように発生するすべての解析結果の尤度値の加算は、シーケンス「one six one e i g h t s i x t h r e e t w o」の尤度であり、これは次式に等しい。

【0065】

【数30】

19

(k)

$$\alpha(6, 3) \times p(o_n_e | s_i_x) \times \beta(9, 3)$$

【0066】ここで、前方向の変数 α の定義により、変数 $\alpha(6, 3)$ はシーケンス「one_s_i_x」の尤度であり、後ろ方向の変数 β の定義により、変数 $\beta(9, 3)$ は、シーケンス「o_n_e」が与えられたときの、シーケンス「e_i_g_h_t_s_i_x_t_h_r_e_e_t_w_o」の尤度である。
(k)

$$\alpha(17, 3) \times p(t_h_r_e_e | s_i_x) \times \beta(22, 5)$$

【0069】ここで、前方向の変数 α の定義により、変数 $\alpha(17, 3)$ はシーケンス「one_s_i_x」の尤度であり、後ろ方向の変数 β の定義により、変数 $\beta(22, 5)$ は、シーケンス「t_h_r_e_e」が与えられたときの、シーケンス「t_w_o」の尤度である。

【0070】従って、「one_s_i_x_t_h_r_e_e_t_w_o」なる学習データにおける、反復回数パラメータ(k+1)における頻度確率 $p(o_n_e | s_i_x)$ に対する、フォワード・バックワードアルゴリズムを用いた再推定式は次式のようにになる。

$$p_i$$

(k)

$$= \alpha(6, 3) \times p(o_n_e | s_i_x) \times \beta(9, 3)$$

(k)

$$+ \alpha(17, 3) \times p(t_h_r_e_e | s_i_x) \times \beta(22, 5)$$

【0072】以上説明したように、本発明の実施形態における特徴は、フォワード・バックワードアルゴリズムを用いて、数23及び数24を含む数22を定式化したことにあるが、当該特徴とする数式は、以下の意味を有する。当該式は、入力データにおいて、当該単位列である第2の単位列が第1の単位列に続くときの単位列のシーケンス間のバイグラムの頻度確率を、上記入力データにおける処理対象の各単位列に対して計算するための式であり、上記シーケンス間のバイグラムの頻度確率は、第1と第2の単位列を含むすべてのセグメント化での尤度の和を、第1の単位列を含むすべてのセグメント化での尤度の和で除算することによって得られる。また、上記式は、上記入力データにおいて各単位列が発生する平均回数を示す分母と、上記入力データにおいて第2の単位列が第1の単位列に続くときの各単位列に対する平均回数を示す分子とを有し、上記分子は、処理対象の各単位列に対する、上記前方尤度と、当該単位列の直前の単位列を条件としたときの当該単位列の頻度確率と、上記後方尤度の積の和であり、上記分母は、処理対象の各単位列に対する、上記前方尤度と、当該単位列の直前の単位列を条件としたときのすべての単位列の頻度確率と、上記後方尤度の積の和である。

20

*る。

【0067】次いで、時刻インデックス(17)においてシーケンス「s_i_x」が終了するすべての解析結果の尤度値の加算は、シーケンス「one_s_i_x_t_h_r_e_e_t_w_o」の尤度であり、これは次式に等しい。

【0068】

【数31】

※【0071】

【数32】

(k+1)

$$p(o_n_e | s_i_x) = p_e / p_i$$

ここで、

【数33】

(k)

$$p_e = \alpha(6, 3) \times p(o_n_e | s_i_x) \times \beta(9, 3)$$

※【数34】

【0073】<<A1. 3>>前方向の変数 α と後ろ方向の変数 β の計算例

例として、データ「one_s_i_x_t_h_r_e_e_t_w_o」について変数 $\alpha(9, 3)$ と変数 $\beta(9, 3)$ を以下に計算する。ここで、変数 $\alpha(9, 3)$ は、シーケンス「one_s_i_x」の尤度であり、このシーケンスは、時刻インデックス9までのシーケンスであって、最後尾において長さ3のシーケンスを有する。また、変数 $\beta(9, 3)$ は、シーケンス「o_n_e」が与えられたときの、シーケンス「e_i_g_h_t_s_i_x_t_h_r_e_e_t_w_o」の条件つき尤度であり、このシーケンスは、時刻インデックス9以降のシーケンスであって、先行するシーケンス「o_n_e」は予め知られている。

【0074】シーケンス「o_n_e」までの尤度(前方の変数) $\alpha(9, 3)$ は、次式で計算される。なお、シーケンス(系列)の長さの最大値を「5」に指定した場合について考える。

【数35】 $\alpha(9, 3)$ =下記の加算値

(a) $n_e_s_i_x$ について: $\alpha(6, 5) \times p(o_n_e | n_e_s_i_x)$

50 (b) $e_s_i_x$ について: $\alpha(6, 4) \times p(o_n_e | e_s_i_x)$

s_i_x)

(c) s_i_xについて: $\alpha(6, 3) \times p(o_n_e | s_i_x)$

(d) i_xについて: $\alpha(6, 2) \times p(o_n_e | i_x)$

(e) xについて: $\alpha(6, 1) \times p(o_n_e | x)$

【0075】シーケンス”o_n_e”の条件のもとでのその後方の尤度(後方の変数) $\beta(9, 3)$ は、次式で計算される。

【数36】 $\beta(9, 3)$ =下記の加算値

(a) e_i_g_h_tについて: $p(e_i_g_h_t | o_n_e) \times$ 10

$\beta(9+5, 5)$ (b) e_i_g_hについて: $p(e_i_g_h$

$| o_n_e) \times \beta(9+4, 4)$ (c) e_i_gについて: p

$(e_i_g | o_n_e) \times \beta(9+3, 3)$ (d) e_iについ

て: $p(e_i | o_n_e) \times \beta(9+2, 2)$

(e) eについて: $p(e | o_n_e) \times \beta(9+1, 1)$

【0076】<<A2>>クラスの事例

シーケンスがクラスに属するケースでは、上述の例のバイグラムの確率部分を、以下のように置き換えることによって変数 α, β が計算される。

(a) $p(o_n_e | n_e_s_i_x)$ は、 $p(\text{class of } o_n_e$ 20
 $| \text{class of } n_e_s_i_x) \times p(o_n_e | \text{class of } o_n_e)$
と取って換えられる。

(b) $p(o_n_e | e_s_i_x)$ は、 $p(\text{class of } o_n_e$
 $| \text{class of } e_s_i_x) \times p(o_n_e | \text{class of } o_n_e)$
と取って換えられる。(c) $p(o_n_e | s_i_x)$

は、 $p(\text{class of } o_n_e | \text{class of } s_i_x) \times p(o_n_e$
 $| \text{class of } o_n_e)$ と取って換えられる。

(d) $p(o_n_e | i_x)$ は、 $p(\text{class of } o_n_e | \text{clas}$
 $s \text{ of } i_x) \times p(o_n_e | \text{class of } o_n_e)$ と取って換

えられる。
(e) $p(o_n_e | x)$ は、 $p(\text{class of } o_n_e | \text{class}$
 $\text{of } x) \times p(o_n_e | \text{class of } o_n_e)$ と取って換えら

れる。
(f) $p(e_i_g_h_t | o_n_e)$ は、 $p(\text{class of } e_i_g$
 $_h_t | \text{class of } o_n_e) \times p(e_i_g_h_t | \text{class of } e_i$

$_g_h_t)$ と取って換えられる。
(g) $p(e_i_g_h | o_n_e)$ は、 $p(\text{class of } e_i_g_h$
 $| \text{class of } o_n_e) \times p(e_i_g_h | \text{class of } e_i_g$

$_h)$ と取って換えられる。
(h) $p(e_i_g | o_n_e)$ は、 $p(\text{class of } e_i_g | \text{cl}$ 40
 $\text{ass of } o_n_e) \times p(e_i_g | \text{class of } e_i_g)$ と取

って換えられる。
(i) $p(e_i | o_n_e)$ は、 $p(\text{class of } e_i | \text{class}$
 $\text{of } o_n_e) \times p(e_i | \text{class of } e_i)$ と取って換えら

れる。
(j) $p(e | o_n_e)$ は、 $p(\text{class of } e | \text{class of } o$
 $_n_e) \times p(e | \text{class of } e)$ と取って換えられる。

【0077】<統計的言語モデル生成処理>図3は、図1の統計的言語モデル生成部20によって実行される統計的言語モデル生成処理を示すフローチャートである。 50

ここで、統計的言語モデル生成部20は、図1に示すように、次のメモリ31乃至36に区分されたワーキングRAM30を備える。

(a) パラメータメモリ31: 当該生成処理で用いる種々の設定パラメータを記憶するメモリである。

(b) シーケンス頻度確率メモリ32: 計算された各シーケンスの頻度確率を記憶するメモリである。

(c) クラス定義メモリ33: 推定された各クラスに属する文字列を記憶するメモリである。

(d) クラス条件付き頻度確率メモリ34: 推定された各クラスに属する各文字列に対する頻度確率、すなわち、クラスの条件付きのクラス間の文字列の頻度確率を記憶するメモリである。

(e) クラスバイグラム頻度確率メモリ35: クラスのバイグラムの頻度確率を記憶するメモリである。

(f) セグメント化されたシーケンスメモリ36: 再推定処理後のセグメント化されたシーケンス(文字列)を記憶するメモリである。

【0078】図3において、まず、ステップS1では、学習用テキストデータメモリ21からテキストデータを読み込む。ここで、入力される学習用テキストデータは、離散的な単位のシーケンスであり、ここで、単位とは例えば、文字であり、シーケンスは単語又は文となり得る文字列である。また、予め下記の入力パラメータが設定されてパラメータメモリ31に記憶されている。

(a) シーケンスの最大長(単位の数で表す。)、

(b) 再推定処理後のクラス数、(c) 廃棄するシーケンス数のしきい値(すなわち、廃棄するシーケンスの発生数の最小値)、及び(d) 終了条件。ここで、終了条件は、例えば、反復回数kのしきい値である。

【0079】次いで、ステップS2で、初期化処理が実行される。入力された学習用テキストデータにおいて、複数の単位からなるシーケンスの相対的な頻度を計数して、それに基づいて各シーケンスの頻度確率を初期設定する。また、上記設定された廃棄するシーケンス数のしきい値以下のシーケンスについては廃棄する。そして、反復回数パラメータkを0にリセットする。

【0080】次いで、ステップS3では、ブラウンアルゴリズムを用いた分類処理を実行する。この分類処理では、反復回数パラメータkのときの各シーケンスの頻度確率に基づいて、クラス間の相互情報量の損失が最小となるように、反復回数パラメータkのときの、クラス定義、クラス条件付きクラス間のシーケンスの頻度確率、及びクラスバイグラムの頻度確率を計算してそれぞれメモリ32乃至35に出力して記憶する。この処理における分類基準は、隣接するシーケンス間の相互情報量であり、上述のアルゴリズムを用いる。これらの相互情報量とアルゴリズムは、隣接する単語の場合に対して、ブラウンによって提案されており、本実施形態では、ブラウンアルゴリズムを用いる。しかしながら、本発明はこれ

に限らず、単位の頻度確率を基礎とする他の分類アルゴリズムを使用することができる。

【0081】次いで、ステップS4において、フォワード・バックワードアルゴリズムを参照して得られた数22-数24を用いて、バイーマルチグラムを用いた再推定処理を実行する。この処理では、直前のステップS3で計算された、反復回数パラメータkのときの、クラス定義、クラス条件付きクラス間のシーケンスの頻度確率、及びクラスバイグラムの頻度確率に基づいて、次の反復パラメータのときのシーケンス間のバイグラムの頻度確率の最尤推定値を得るように、反復回数パラメータ(k+1)のときの、各シーケンスの頻度確率を再推定して計算して、メモリ32に出力して記憶する。この処理における処理基準は、上記数22-数24を用いて、すなわち、複数のシーケンスのクラスとバイグラムの依存性を仮定して計算された解析結果の尤度の中の最大値である最尤推定値を基準値として用いることであり、再推定のためのアルゴリズムとしてEMアルゴリズムを用いる。

【0082】次いで、ステップS5で、所定の終了条件を満足するか否かが判断され、NOのときは、ステップS6で反復回数パラメータkを1だけインクリメントしてステップS3及びS4の処理を繰り返す。一方、ステップS5でYESであれば、生成された統計的言語モデルのデータを統計的言語モデルメモリ22に出力して記憶する。ここで、生成された統計的言語モデルのデータとは、各シーケンスの頻度確率に関するデータであり、具体的には、下記のデータである。

(a) 入力されたデータを複数のシーケンスにセグメント化したときの最尤推定値を有する各シーケンスのデータ；

(b) クラス定義、すなわち、各クラスにおけるシーケンス；及び

(c) クラスの頻度確率、すなわち、各クラスのバイグラム確率、各シーケンスのクラス条件付き確率。

【0083】図4は、図3のサブルーチンであるブラウンアルゴリズムを用いた分類処理を示すフローチャートである。単語の自動分類のために、ブラウン他によってシーケンスの自動分類に使用するためのアルゴリズム

(例えば、従来技術文献5参照。)が提案されており、本実施形態では、これを使用する。ブラウンらは、文章の尤度を最大化するクラスへの分割又はセグメント化

が、隣接する単語間の相互情報量を最大化する分割又はセグメント化でもあることを示している。彼らは単語のバイグラム分布を入力とし、単語クラスへの分割及びクラス分布を出力する貪欲なアルゴリズム(greedy algorithm)を提案している。一方、本発明者は、入力としてバイーマルチグラムの頻度確率の分布(すなわち、シーケンスのバイグラムの頻度確率の分布)を採用することにより、このアルゴリズムを適用している。出力は、シーケンスのクラスへのセグメント化及びその各シーケンスの頻度確率の分布である。

【0084】この分類処理で用いる相互情報量を用いた単語のクラスタリングについて詳細説明する(例えば、従来技術文献6「北研二ほか著、"音声言語処理", 森北出版, pp. 110-113, 1996年11月15日発行」参照。)。ここでは、隣接する単語に基づく単語の分類法として、クラス間の相互情報量を最大にする方法について説明する。相互情報量に基づくクラスタリングは、バイグラムのクラスモデルにおいて単語をクラスへ分割する最尤な方法は、隣接するクラスの平均相互情報量を最大にするようなクラス割り当てであることを、理論的な根拠としている。N-gramのクラスモデルとは、次式のように、単語のクラスのN-gramとクラス別の単語の出現分布の組み合わせで、単語のN-gramを近似する言語モデルのことである(この式は、単語クラスを品詞に置き換えれば、形態素解析におけるHMMの式と同じになる。従って、この単語分類法は、最適な品詞体系を自動的に求める方法とも考えられる。

【数37】 $P(w_i | w_1^{i-1}) \doteq P(w_i | c_i) P(c_i | c_1^{i-1})$

【0085】ここで、単語 w_i をクラス c_i に写像する関数 π を用いて、V個の単語をC個のクラスに分割すると仮定する。学習テキスト t_1^T が与えられたとき、 $P(t_2^T | t_1) = P(T_2 | T_1) P(t_3 | t_2) \cdots P(t_T | t_{T-1})$ を最大にするように関数 π を決めればよい。詳細は省略するが、単語あたりの対数尤度 $L(\pi)$ 、単語のエントロピー $H(w)$ 、隣接するクラスの平均相互情報量 $I(c_1; c_2)$ の間には、近似的に次式の関係が成り立つ。

【0086】

【数38】

$$\begin{aligned}
& L(\pi) \\
&= (T-1)^{-1} \log P(t_2^T | t_1) \\
&= \sum_{w_1, w_2} \{C(w_1 w_2)\} / (T-1) \times \log P(c_2 | c_1) P(w_2 | c_2) \\
&\quad \times P(c_1 c_2) \log \{P(c_2 | c_1) / P(c_2)\} \\
&\quad + \sum_w P(w) \log P(w) \\
&= I(c_1; c_2) - H(w)
\end{aligned}$$

【0087】ここで、 $H(w)$ は分割 π に依存しないから、 $L(\pi)$ を最大化するためには、 $I(c_1; c_2)$ を最大化すればよい。いまのところ、平均相互情報量を最大化するような分割を求めるアルゴリズムは知られていない。しかしながら、本実施形態で用いる次のような貪欲なアルゴリズム (greedy algorithm) でも、かなり興味深いクラスタを得ることができる。このように包含関係を持つクラスタを生成する方法は、階層的クラスタリングと呼ばれる。これに対して、 k 平均アルゴリズムのように、重なりを持たないクラスタを生成する方法は非階層的クラスタリングと呼ばれる。

【0088】次の併合を $V-1$ 回繰り返すと、すべての単語が一つのクラスになる。すなわち、クラスが併合される順序から、単語を葉とする二分木ができる。

1. すべての単語に対して、一つのクラスを割り当てる。
2. 可能な二つのクラスの組み合わせの中で、平均相互情報量の損失を最小にする組み合わせを選択し、これらを一つのクラスに併合する。
3. ステップ2を $V-C$ 回繰り返すと C 個のクラスが得られる。

【0089】一般に、クラスタが形成される過程を表す階層構造は樹形図 (dendrogram) と呼ばれるが、自然言語処理ではこれをソーラスの代わりに使うことができる。単純に考えると、この準最適アルゴリズムは、語彙数 V に対して V^3 の計算量を必要とする。しかし、

(1) 二つのクラスタを併合したときの情報量の変化だけを求めればよいことや、(2) 二つのクラスタの併合により相互情報量が増加するのは全体の一部に過ぎないことを利用すれば、 $O(V^3)$ の計算、すなわち、繰り返し回数 V の三乗に比例するオーダーの計算コストで済む。

【0090】分類処理 (又はクラスタリング処理) を示す図4において、まず、ステップS11では、初期設定処理が実行され、各シーケンスをその自らのクラスに割り当てる。すなわち、各シーケンス s_i それぞれ各クラス C_i に割り当てる。従って、クラスの初期バイグラム の頻度確率の分布はシーケンスのバイグラムの頻度確率の分布に等しく、また、

$$p(s_i | C_i) = 1$$

である。

【0091】次いで、ステップS12で、各クラスの対 (C_k, C_l) について、クラス C_k とクラス C_l とをマージしたときの相互情報量の損失を計算した後、ステップS13で、相互情報量の損失が最小であるクラスの対をマージする。そして、ステップS14で、上記マージに従って、メモリ34及び35に記憶されたクラスの頻度確率の分布を更新する。次いで、ステップS15で、ステップS2の初期化処理で設定された必要なクラス数が得られたか否かが判断され、NOであるときは、ステップS12に戻り、上記の処理を繰り返す。一方、ステップS15で、YESのときは、元のメインルーチンに戻る。

【0092】<音声認識装置>次いで、図1に示す連続音声認識装置の構成及び動作について説明する。図1において、単語照合部4に接続された音素隠れマルコフモデル (以下、隠れマルコフモデルをHMMという。) メモリ11内の音素HMMは、各状態を含んで表され、各状態はそれぞれ以下の情報を有する。

(a) 状態番号、(b) 受理可能なコンテキストクラス、(c) 先行状態、及び後続状態のリスト、(d) 出力確率密度分布のパラメータ、及び(e) 自己遷移確率及び後続状態への遷移確率。なお、本実施形態において用いる音素HMMは、各分布がどの話者に由来するかを特定するため、所定の話者混合HMMを変換して生成する。ここで、出力確率密度関数は3次元の対角共分散行列をもつ混合ガウス分布である。また、単語照合部4に接続された単語辞書メモリ12内の単語辞書は、音素HMMメモリ11内の音素HMMの各単語毎にシンボルで表した読みを示すシンボル列を格納する。

【0093】図1において、話者の発声音声はマイクロホン1に入力されて音声信号に変換された後、特徴抽出部2に入力される。特徴抽出部2は、入力された音声信号をA/D変換した後、例えばLPC分析を実行し、対数パワー、16次ケプストラム係数、 Δ 対数パワー及び16次 Δ ケプストラム係数を含む3次元の特徴パラメータを抽出する。抽出された特徴パラメータの時系列はバッファメモリ3を介して単語照合部4に入力される。

【0094】単語照合部4は、ワンパス・ビタビ復号化法を用いて、バッファメモリ3を介して入力される特

徴パラメータのデータに基づいて、音素HMM11と単語辞書12とを用いて単語仮説を検出し尤度を計算して出力する。ここで、単語照合部4は、各時刻の各HMMの状態毎に、単語内の尤度と発声開始からの尤度を計算する。尤度は、単語の識別番号、単語の開始時刻、先行単語の違い毎に個別にもつ。また、計算処理量の削減のために、音素HMM11及び単語辞書12とに基づいて計算される総尤度のうちの低い尤度のグリッド仮説を削減する。単語照合部4は、その結果の単語仮説と尤度の情報を発声開始時刻からの時間情報（具体的には、例えばフレーム番号）とともにバッファメモリ5を介して単語仮説絞込部6に出力する。

【0095】単語仮説絞込部6は、単語照合部4からバッファメモリ5を介して出力される単語仮説に基づいて、統計的言語モデルメモリ22内の統計的言語モデルを参照して、終了時刻が等しく開始時刻が異なる同一の単語の単語仮説に対して、当該単語の先頭音素環境毎に、発声開始時刻から当該単語の終了時刻に至る計算された総尤度のうちの最も高い尤度を有する1つの単語仮説で代表させるように単語仮説の絞り込みを行った後、絞り込み後のすべての単語仮説の単語列のうち、最大の総尤度を有する仮説の単語列を認識結果として出力する。本実施形態においては、好ましくは、処理すべき当該単語の先頭音素環境とは、当該単語より先行する単語仮説の最終音素と、当該単語の単語仮説の最初の2つの音素とを含む3つの音素並びをいう。

【0096】例えば、図2に示すように、 $(i-1)$ 番目の単語 W_{i-1} の次に、音素列 a_1, a_2, \dots, a_n からなる i 番目の単語 W_i がくるときに、単語 W_{i-1} の単語仮説として6つの仮説 $W_a, W_b, W_c, W_d, W_e, W_f$ が存在している。ここで、前者3つの単語仮説 W_a, W_b, W_c の最終音素は $/x/$ であるとし、後者3つの単語仮説 W_d, W_e, W_f の最終音素は $/y/$ であるとする。終了時刻 t_e と先頭音素環境が等しい仮説（図2では先頭音素環境が $"x/a_1/a_2"$ である上から3つの単語仮説）のうち総尤度が最も高い仮説（例えば、図2において1番上の仮説）以外を削除する。なお、上から4番めの仮説は先頭音素環境が違うため、すなわち、先行する単語仮説の最終音素が x ではなく y であるので、上から4番めの仮説を削除しない。すなわち、先行する単語仮説の最終音素毎に1つのみ仮説を残す。図2の例では、最終音素 $/x/$ に対して1つの仮説を残し、最終音素 $/y/$ に対して1つの仮説を残す。

【0097】以上の実施形態においては、当該単語の先*

* 頭音素環境とは、当該単語より先行する単語仮説の最終音素と、当該単語の単語仮説の最初の2つの音素とを含む3つの音素並びとして定義されているが、本発明はこれに限らず、先行する単語仮説の最終音素と、最終音素と連続する先行する単語仮説の少なくとも1つの音素とを含む先行単語仮説の音素列と、当該単語の単語仮説の最初の音素を含む音素列とを含む音素並びとしてもよい。

【0098】以上の実施形態において、特徴抽出部2と、単語照合部4と、単語仮説絞込部6と、統計的言語モデル生成部20とは、例えば、デジタル電子計算機などのコンピュータで構成され、バッファメモリ3、5と、音素HMMメモリ11と、単語辞書メモリ12と、学習用テキストデータメモリ21と、統計的言語モデルメモリ22とは、例えばハードディスクメモリなどの記憶装置で構成される。

【0099】以上実施形態においては、単語照合部4と単語仮説絞込部6とを用いて音声認識を行っているが、本発明はこれに限らず、例えば、音素HMM11を参照する音素照合部と、例えばOne Pass DPアルゴリズムを用いて統計的言語モデルを参照して単語の音声認識を行う音声認識部とで構成してもよい。

【0100】

【実施例】<統計的言語モデル生成処理の第1の実施例>入力される学習データが、以下のような1000文字列の場合であり、単位である文字から単語にセグメント化するための例である。

「onesixoneeightfivezero
...」

但し、奇数の単語の後には必ず偶数の単語が後続し、偶数の単語の後には必ず奇数の単語が後続する場合である。当該実施例における入力パラメータは以下の通りである。

(a) 1個のシーケンスの最大長=5、(b) クラス数=2、及び(c) 廃棄するシーケンスのしきい値=100。

【0101】初期化処理($k=0$)では、学習データにおいて、100回を越えて観測した文字のすべての組合せの相対的な計数値を初期値とする。従って、反復パラメータ $k=0$ におけるシーケンスの頻度確率の分布の計数結果は次の表ようになる。なお、各シーケンスの $n_b(\cdot)$ は計数値を表す。

【0102】

【表2】

$$\begin{aligned} p(n|o) &= n_b(on) / n_b(o) = 0.08 \\ p(n_e|o) &= n_b(one) / n_b(o) = 0.06 \\ \dots \\ p(n_{e_s_i_x}|o) &= n_b(onesix) / n_b(o) = 0.005 \end{aligned}$$

$p(e|o_n) = nb(one) / nb(on) = 0.9$
 $p(e_s|o_n) = nb(ones) / nb(on) = 0.005$
 \dots
 $p(e_s_i_x_o|o_n) = nb(onesix) / nb(on) = 0.001$
 \dots
 $p(s_i_x|o_n_e) = nb(onesix) / nb(one) = 0.05$
 \dots

【0103】ステップS3の分類処理では、入力データは、反復パラメータ $k=0$ のときのシーケンスの頻度確率の分布であり、当該分類処理における出力データは、以下のようになる。

(a) 反復パラメータ $k=1$ のときのクラス定義

【数40】 $class1 = \{e_s_i_x_o; e_t_w_o; n_e_s_i_x; \dots; f_o_u_r; f_o_u_r_f; \dots; g_h_t_s; g_h_t_o_n_e; e_i_g_h_t\}$

【数41】 $class2 = \{o_n_e; e_s_i_x_o; x; f_i_v_e; f_i_v_e; t_s_e_v; s_e_v_e_n; \dots; x_n_i; x_n_i_n_e; n_i_n_e; \dots\}$

$class3 = \dots$

(b) 反復パラメータ $k=1$ のときのクラス条件付き頻度確率の分布

【数42】

$p(e_s_i_x_o|class1), p(e|class1), \dots$

$p(o_n_e|class2), p(e_s_i_x_o|class2), \dots$

(c) 反復パラメータ $k=1$ のときのクラスバイグラムの頻度確率の分布

【数43】 $p(class1|class2) = 0.3$

$p(class2|class1) = 0.1$

$p(class3|class1) = 0.4$

\dots

【0104】ステップS4の再推定処理では、反復パラメータ $k=1$ のときのクラス定義及びクラスの頻度確率の分布を入力データとし、次に示す反復パラメータ $k=1$ のときのシーケンスの頻度確率の分布を出力する。

【数44】 $p(n|o) = 0.9$

$p(n_e|o) = 0.8$

$p(n_e_s|o) = 0.05$

\dots

$p(n_e_s_i_x|o) = 0$

【数45】 $p(e|o_n) = 0.02$

$p(e_s|o_n) = 0.001$

\dots

$p(e_s_i_x_o|o_n) = 0$

\dots

$p(s_i_x|o_n_e) = 0.5$

\dots

【0105】以下同様に処理が実行され、第1の実施例

における出力結果は以下のようになる。

(a) セグメント化された入力文字列 (MLセグメント化)

"o_n_e s_i_x o_n_e e_i_g_h_t f_i_v_e z_e_r_o ..."

(b) クラス定義

【数46】 $class1 = \{o_n_e; t_h_r_e_e; f_i_v_e; s_e_v_e_n; n_i_n_e\}$

$class2 = \{z_e_r_o; t_w_o; f_o_u_r; s_i_x; e_i_g_h_t\}$

(c) クラス条件付きの頻度確率の分布

20 【数47】 $p(o_n_e|class1) = 0.2$

$p(t_h_r_e_e|class1) = 0.2$

$p(f_i_v_e|class1) = 0.2$

\dots

$p(z_e_r_o|class2) = 0.2$

$p(t_w_o|class2) = 0.2$

(d) クラスバイグラムの頻度確率の分布

【数48】 $p(class1|class2) = 1$

$p(class2|class1) = 1$

30 【0106】<統計的言語モデル生成処理の第2の実施例>入力される学習データが、自然言語のテキストデータによる以下の文、すなわち単語列である場合であって、単位である単語をフレーズにセグメント化する場合を説明するための実施例である。ここで、<s>は開始を示す記号であり、</s>は終了を示す記号である。

「<s> good afternoon new washington hotel may i help you ...</s>」

ここで、入力パラメータは、以下の通りである。

(a) シーケンスの最大長=数個の単語 (例えば、1乃至5個の単語、以下の実施例では、4)、(b) クラス数=1000、及び (c) 初期化処理のしきい値=3

0。

【0107】初期化処理 ($k=0$) では、学習データにおいて、30回を越えて観測した単語のすべての組合せの相対的な計数値を初期値とする。従って、反復パラメータ $k=0$ におけるシーケンスの頻度確率の分布の計数結果は次の表のようになる。

【0108】

【表3】

```

p (afternoon | good)
=nb (good afternoon) / nb (good) = 0.08
p (afternoon_new | good)
=nb (good afternoon new) / nb (good) = 0.06
p (good_afternoon | <s>)
=nb (<s>good afternoon) / nb (<s>) = 0.06
...
p (</s> | may_i_help_you)
=nb (may i help you </s>) / nb (may i help you)
=0.005

```

【0109】そして、第2の実施例における出力結果は以下のようになる。

(a) セグメント化された入力文字列 (MLセグメント化)

「good_afternoon new_washington_hotel may_i_help_you」

(b) クラス定義

【数49】 $class1 = \{good_afternoon; good_morning; hello; may_i_help_you...\}$

...

$class2 = \{new_washington_hotel; sheraton_hotel; plaza;...\}$

...

$class1000 = \{give_me_some; tell_me\}$

(c) クラス条件付き頻度確率の分布

【数50】

$p(good_afternoon | class1) = 0.003$

$p(good_morning | class1) = 0.002$

$p(hello | class1) = 0.002$

...

(d) クラスバイグラムの頻度確率の分布

【数51】 $p(class2 | class1) = 0.04$

$p(class3 | class1) = 0.005$

...

【0110】＜実験及び実験結果＞本発明者は、実施形態の装置の性能を実験するために、下記の実験を行った。まず、プロトコル及びデータベースの実験及び実験結果について述べる。可変長フレーズ間のバイグラム依存を学習する目的は、従来のワードバイグラムモデルの限界を改善する一方で、モデル内のパラメータ数を単語のトライグラムの場合よりも少なくすることにある。従って、バイマルチグラムモデルの評価を行うために適する基準は、その予測能力、パラメータ数を測定し、従来のバイグラム、トライグラムモデルのそれらと比較することである。予測能力は通常、次式のパープレキシテ*

*イの測定によって評価される。

【0111】

【数52】

$PP = \exp \{-(1/T) \log(L(W))\}$

【0112】ここで、Tを文Wにおける単語の数である。パープレキシティPPが低いほど、モデルの予測がより高精度であることを示す。統計的モデルでは、実際には2つのパープレキシティ値PP及びPP*が存在し、数52の中のL(W)をそれぞれ次式として計算される。

【0113】

【数53】

$$L(W) = \sum_S L(W, S)$$

及び

30 【数54】 $L(W) = L(W, S^*)$

【0114】2つのパープレキシティPP* - PPの差は、常に正の数又は零であり、文Wの解析結果Sの曖昧さの度合い、あるいは発話認識機のように最良の解析結果の尤度を用いて文の尤度に到達する場合は、予測の正確さにおける損失を測定する。

【0115】以下では、先ず、ある推定手順における損失(PP* - PP)を評価し、この推定手順自体の影響力についてフォワード・バックワードアルゴリズム(数15)又は決定論的方法(数16)を用いて考察する。最後に、これら結果を従来のn-gramモデルを用いて得られた結果と比較する。本目的の達成のため、クラークソン (Clarkson) ほか1997年) による公知のCMUツールキットを用いる。実験対象として、次の表の本特許出願人が所有する「旅行の手配」に関するデータを使用する。

【0116】

【表4】

本特許出願人が所有する「旅行の手配」に関するデータ

学習

テスト

文の数	13650	2430
トークンの数	167000	29000 (1%OOV)
語彙数	3525	+2800OV

(注) OOVは、Out Of Vocabularyの略であり、語彙にない単語をいう。

【0117】本データベースは、ホテルのクラークと顧客の間で自発的に行われた旅行／宿泊施設情報についての対話である。言いよどみの単語、及び間違った開始は、単一のマーカー “uh” にマッピングされる。本実験において、フレーズの最大長は $n=1$ 語から4語まで変化させた($n=1$ ではバイマルチグラムは従来のバイグラムに相当する)。すべてのバイマルチグラムの頻度確率は、6回のトレーニング反復で推定され、初期化において20回以下、各反復において10回以下の頻度でしか現れないすべての文を放棄し、フレーズ辞書の枝刈りを行った。ここで、初期化におけるしきい値が10-30の範囲にあるとき、本データにおいて、異なる枝刈り限界値を用いても結果に重大な影響が及ぶことはない。反復の場合のしきい値はその約半分である。

【0118】しかしながら、すべての1単語フレーズは、その推定出現回数にかかわらず維持されるため(フレーズ s_i 及び s_j が1単語フレーズであり、組み合わせ $c(s_i, s_j)$ の再推定値が零であると、組み合わせ $c(s_i, s_j)$ は1にリセットされる。)、すべてのワードバイグラムが最終辞書に現れることになる。さらに、す*非決定性の方式の度合い

* 全ての n -gram及びフレーズのバイグラム確率は、ウィッテン (Witten) ほか (1991年) による公知のWitten-Bellディスカунティング法を用いて、カツ (Katz) (1987年) による公知のバックオフ・スムージング法で平滑化される。ここで、Witten-Bellディスカунティング法を選択したのは、本テストデータにおいて従来の n -gramを用いた場合、最良のパープレキシティスコアが得られるためである。

【0119】次いで、クラスタリングを行わない実験について述べる。まず、非決定性の方式の度合いにおいては、表4の本特許出願人が所有する「旅行の手配」に関するデータに対するテストで、フォワード・バックワードアルゴリズムによる学習の後に得られたパープレキシティ値 PP^* 及び PP を次の表に示す。パープレキシティ値の差($PP^* - PP$)は通常、パープレキシティの約1ポイント以内にとどまる。すなわち、単一の最良フレーズに依存しても、予測の正確さが大幅に損なわれることがあつてはならないことを意味している。

【0120】

【表5】

n	1	2	3	4
PP	56.0	43.9	44.2	45.0
PP*	56.0	45.1	45.4	46.3

【0121】次いで、再推定手順の影響力では、フォワード・バックワードアルゴリズム又はビタビ推定アルゴリズムのいずれかを用いたパープレキシティ値 PP^* 及び推定方法の影響：テストパープレキシティ値 PP^*

※びモデルサイズを次の表に示す。

【0122】

【表6】

n	1	2	3	4
FB法	56.0	45.1	45.4	46.3
ビタビ法	56.0	45.7	45.9	46.2

【0123】

★ ★ 【表7】

推定方法の影響：モデルのサイズ

n	1	2	3	4
---	---	---	---	---

	35			36
FB法	3 2 5 0 5	4 4 3 8 2	4 3 6 7 2	4 3 1 8 6
ビタビ法	3 2 5 0 5	6 5 1 4 1	6 7 2 5 8	6 7 2 9 5

【0124】表6及び表7から明らかなように、パープレキシティ値に関する限り、推定方法はほとんど影響を及ぼさず、フォワード・バックワードアルゴリズムによる学習を用いる方がわずかながら有利であるように見える。一方、モデルのサイズは、学習終了時に個々のバイーマルチグラム数として測定された場合、フォワード・バックワードアルゴリズムによる学習において約30%も減少する。すなわち、同じテストパープレキシティ値に対して、おおよそ40, 000対60, 000の違いとなる。

【0125】バイーマルチグラム結果は、概して、フレーズ放棄を行う枝刈りのための発見的知識では完全に過学習を回避できないことを示唆する。確かに、(おそら*

n-gramの比較

テストパープレキシティ値 P P				
nの値	1	2	3	4
n-gram	3 1 4. 2	5 6. 0	4 0. 4	3 9. 8
バイーマルチグラム	5 6. 0	4 3. 9	4 4. 2	4 5. 0

【0128】

※ ※ 【表9】

n-gramの比較

モデルのサイズ				
n値	1	2	3	4
n-gram	3 5 2 6	3 2 5 0 5	7 5 5 1 1	1 1 2 1 4 8
バイーマルチグラム	3 2 5 0 5	4 4 3 8 2	4 3 6 7 2	4 3 1 8 6

【0129】表8及び表9から明らかなように、最も低いバイーマルチグラムパープレキシティスコア(43.9)は、トライグラムの値よりも依然として高いが、バイグラム値(56.0)よりもトライグラム値(40.4)により近い値となっている。さらに、トライグラムスコアはディスカウントされた方法に依存する。なお、線形ディスカウンティング法では、本テストにおけるトライグラムのパープレキシティは、48.1であった。

【0130】5-グラムのパープレキシティ値(上記表に示さず)は40.8であり、4-gramスコアよりもやや高い。これは、バイーマルチグラムパープレキシティが $n > 2$ (すなわち、依存性が4語以上にわたる場合)のとき減少しないという事実に一致する。最後に、バイーマルチグラムモデルのエントリ数はトライグラムモデルのエントリ数よりも少なく(45000に対して

*くは6から8語にまたがる依存性を意味する) $n = 3$, 4のパープレキシティ値は、(依存性が4語に限定される) $n = 2$ のときのそれよりも高くなる。他の方法、おそらくは短いものよりも長いフレーズを不利にするような方法であれば成功ものと考えられる。

【0126】さらに、n-gramとの比較においては、フォワード・バックワードアルゴリズムによる学習から得られたパープレキシティ値(PP)、n-gramに対するモデルサイズ、及びバイーマルチグラムを次の表に示す。

【0127】

【表8】

75000)、マルチグラムが達成するモデルの正確性とモデルサイズ間のトレードオフが示されている。

【0131】さらに、クラスタリングを用いた実験及び実験結果について述べる。本実験では、フレーズのクラスタリングによってパープレキシティスコアは改善されなかった。パープレキシティの増加が非常に少なくなる(1ポイント以下)のは、フレーズのほんの一部(10~20%)のみがクラスタとなる時であり、これを越えるとパープレキシティはかなり悪化する。この効果は、クラス推定が単語推定に統合されない時、n-gramの枠組みにおいても度々報告されている。しかしながら、フレーズのクラスタリングによって、自然発話の特徴づける言いよどみの語の挿入等、ことばの非流暢性のいくつかを自然に扱うことができる。この点を説明するために、先ず $n = 4$ 語までのフレーズを扱うモデルの学

習の間に統合されるフレーズを次の表に列挙する。こ
で、言いよどみを示す “u h” を含むフレーズはこの
表の上部に示す。主に、話者の言いよどみによるフレーズ

* ズの違いは、共に統合されることが多い。

【0132】

【表10】

4語シーケンスまでを扱うモデルにおける統合されたフレーズの一例

```
{yes_that_will;uh_that_would}
{yes_that_will_be;uh_yes_that's}
{uh_by_the;and_by_the}
{yes_uh_i;i_see_i}
{okay_i_understand;uh_yes_please}
{could_you_recommend;uh_is_there}
{uh_could_you_tell;and_could_you_tell}
{so_that_will;yes_that_will;yes_that_would;uh_that_would}
{if_possible_i'd_like;we_would_like;uh_i_want}
{that_sounds_good;uh_i_understand}
{uh_i_really;uh_i_don't}
{uh_i'm_staying;and_i'm_staying}
{all_right_we;uh_yes_i}
```

```
{good_morning_this;good_afternoon_this}
{yes_i_do;yes_thank_you}
{we'll_be_looking_forward;we_look_forward}
{dollars_a_night;and_forty_yen}
{for_your_help;for_your_information}
{hold_the_line;want_for_a_moment}
{yes_that_will_be;and_could_you_tell}
{please_go_ahead;you_like_to_know}
{want_time_would_you;and_you_would}
{yes_there_is;but_there_is}
{join_phillips_in_room;ms._suzuki_in}
{name_is_suzuki;name_is_ms._suzuki}
{i'm_calling_from;a;also_i'd_like}
{much_does_it_cost;can_reach_you}
{thousand_yen_room;dollars_per_person}
{yes_i_do;yes_thank_you;i_see_sir}
{you_tell_me_where;you_tell_me_what}
{a_reservation_for_the;the_reservation_for}
{your_name_and_the;you_give_me_the}
{amy_harris_in;is_amy_harris_in}
{name_is_mary_phillips;name_is_kazuo_suzuki}
{hold_on_a_moment;wait_a_moment}
{give_me_some;also_tell_me}
```

【0133】カワハラ (Kawahara) ら (1997年) に
よれば、上記の表はさらに、単語予測とは別に、フレー
ズ検索及びクラスタリングを行う他の動機づけ、すなわ
ちトピックの識別や対話のモデリング、及び言語理解に
関する問題への対応を示している。確かに本実験におけ
るクラスタとなったフレーズは、完全盲目的、すなわち
意味論的／語用論的情報を全くなくして導かれたもので

あるが、クラス内フレーズには強固な意味論的相関関係
が示されている。しかしながら、本手法を音声理解に効
率的に使用できるようにするためには、拘束条件は、例
えばスピーチアクトタグ (speech act tags) のような
いくつかのより高いレベルの情報を用いてフレーズクラ
スタリング処理に設定する必要がある。

【0134】以上説明したように、フレーズ間に $n - g$

ram依存を仮定する可変長フレーズを導くアルゴリズムは、言語モデリングのタスクのために提案され、推定されてきた。特定タスクの言語コーパスは、文をフレーズに構成することによりバイグラムパープレキシティ値を大幅に減らし、一方で言語モデルにおけるエントリ数をトライグラムモデルの場合に比べてより低い値に保つことが可能であることを示している。しかしながら、これら結果は、より効率的な枝刈り方法によってさらに改善され、不要な学習を行わずにより長い依存性について学習することが可能となる。さらに、語形変化の態様を簡単に本枠組み内に統合することができるため、異なる長さを有するフレーズに共通のラベルを割り当てることが可能である。フレーズの意味論的關係が統合されるので、本手法は対話モデリングや言語理解の分野においても用いられる。その場合、意味論的／語用論的情報を用いれば、フレーズクラスを得るための処理に制限を設けることができる。

【0135】＜変形例＞以上の実施形態においては、単位は英語の文字であり、シーケンスは単語であり、上記分類処理は、文字列を複数の単語の列に分類し、上記統計的シーケンスモデルは、統計的言語モデルである。本発明はこれに限らず、単位は、日本語などの他の自然言語の文字であってもよい。また、単位は自然言語の単語であり、シーケンスはフレーズであり、上記分類処理は、単語列を複数のフレーズの列に分類し、上記統計的シーケンスモデルは、統計的言語モデルであってもよい。

【0136】＜実施形態の効果＞以上説明したように、本発明に係る実施形態によれば、以下のような特有の効果を有する。

(A) EMアルゴリズムを使用して単語のシーケンスの頻度分布を計算することができ、ML基準を最適化することができる。すなわち、本実施形態のアルゴリズムを用いられれば、必ず、クラスターリングの処理を単調収束させることができ、最適値の解析結果を得ることができる。

(B) シーケンス分類の解析を自由にすることができる。具体的には、上述のフォワード・バックワードアルゴリズムを用いた非決定性の手法を用いるので、自由度のある解が得られる。なお、当該非決定性の手法を用いることができるのは、変数 α 、 β を決めることができるからである。従って、入力データの尤度を改善することにより、シーケンス [b c d] が入力シーケンスにあったときに、[b c] + [d]、[b] + [c d]、

[b] + [c] + [d] 等の小シーケンスへの分割が可能である。言い換えれば、あるシーケンスが入力シーケンスに与えられていても、解析は事前に決定されず、すべては入力データの尤度に依存する、つまり確定的ではなく、入力データの頻度確率に依存してクラスターリングの処理が行われる。

(C) 可変長のシーケンスの自動的分類を行うことができる。ここで、シーケンスの分類を、単語の分類に依存させない。また、シーケンスの分類を直接的に自動的に行なって、長さの違う共通のクラスシーケンスに高精度で分類できる。

【0137】従って、本発明に係る実施形態によれば、従来例に比較して、最適な状態に向かう単調な収束を保証することができ、自由度があり、可変長のシーケンスを同一のクラスで取り扱うことができ、デジタル計算機を用いて実用的に高速処理することができる統計的シーケンスモデル生成装置、統計的言語モデル生成装置及び音声認識装置を提供することができる。

【0138】

【発明の効果】以上詳述したように本発明に係る統計的シーケンスモデル生成装置によれば、1個又は複数の単位からなる単位列であるシーケンスを含む入力データに基づいて、可変長の自然数 N_1 個の単位列であるマルチグラムと、可変長の自然数 N_2 個の単位列であるマルチグラムとの間のバイグラムであるバイーマルチグラムの統計的シーケンスモデルを生成する統計的シーケンスモデル生成装置であって、上記入力データに基づいて、予め決められた N_1 、 N_2 の最大値の拘束条件のもとで、すべての単位列の組み合わせの上記バイグラムの頻度確率を計数する初期化手段と、上記初期化手段によって計数された上記バイグラムの頻度確率に基づいて、各クラスの対をマージしたときの相互情報量の損失が最小となるようにマージして各クラスの頻度確率を更新して予め決められた数の複数のクラスに分類することにより、分類されたクラスに含まれる単位列と、分類されたクラスの条件付きの単位列の頻度確率と、分類されたクラス間のバイグラムの頻度確率を計算して出力する分類手段と、上記分類処理手段から出力される分類されたクラスに含まれる単位列と、分類されたクラスの条件付きの単位列の頻度確率と、分類されたクラス間のバイグラムの頻度確率とに基づいて、EMアルゴリズムを用いて、最尤推定値を得るように再推定し、ここで、フォワード・バックワードアルゴリズムを用いて、処理対象の各単位列に対して、時系列的に前方にとり得る処理対象の当該単位列に対する前方尤度と、当該単位列の直前の単位列を条件としたときの当該単位列の頻度確率と、時系列的に後方にとり得る当該単位列に対する後方尤度とに基づいてシーケンス間のバイグラムの頻度確率を示す式を用いて、当該シーケンス間のバイグラムの頻度確率を再推定することにより、再推定結果である上記バイーマルチグラムの統計的シーケンスモデルを生成して出力する再推定手段と、上記分類手段の処理と上記再推定手段の処理を所定の終了条件を満たすまで繰り返し実行するように制御する制御手段とを備える。従って、本発明によれば、従来例に比較して、最適な状態に向かう単調な収束を保証することができ、自由度があり、可変長のシーケ

ンスを同一のクラスで取り扱うことができ、デジタル計算機を用いて実用的に高速処理して統計的シーケンスモデルを生成することができる統計的シーケンスモデル生成装置を提供することができる。

【0139】また、本発明に係る統計的言語モデル生成装置によれば、上記統計的シーケンスモデル生成装置において、上記単位は自然言語の文字であり、上記シーケンスは単語であり、上記分類手段は、文字列を複数の単語の列に分類し、上記統計的シーケンスモデルは、統計的言語モデルである。従って、本発明によれば、従来例に比較して、最適な状態に向かう単調な収束を保証することができ、自由度があり、可変長のシーケンスを同一のクラスで取り扱うことができ、デジタル計算機を用いて実用的に高速処理して統計的言語モデルを生成することができる統計的言語モデル生成装置を提供することができる。

【0140】さらに、本発明に係る統計的言語モデル生成装置によれば、上記統計的シーケンスモデル生成装置において、上記単位は自然言語の単語であり、上記シーケンスはフレーズであり、上記分類手段は、単語列を複数のフレーズの列に分類し、上記統計的シーケンスモデルは、統計的言語モデルである。従って、本発明によれば、従来例に比較して、最適な状態に向かう単調な収束を保証することができ、自由度があり、可変長のシーケンスを同一のクラスで取り扱うことができ、デジタル計算機を用いて実用的に高速処理して統計的言語モデルを生成することができる統計的言語モデル生成装置を提供することができる。

【0141】またさらに、本発明に係る音声認識装置によれば、入力される発声音声文の音声信号に基づいて、所定の統計的言語モデルを用いて音声認識する音声認識手段を備えた音声認識装置において、上記音声認識手段は、上記統計的言語モデル生成装置によって生成された統計的言語モデルを参照して音声認識する。従って、本発明によれば、従来例に比較して、最適な状態に向かう

単調な収束を保証することができ、自由度があり、可変長のシーケンスを同一のクラスで取り扱うことができ、デジタル計算機を用いて実用的に高速処理して統計的言語モデルを生成することができる。また、当該生成された統計的言語モデルを用いて音声認識することにより、従来例に比較して高い音声認識率で音声認識することができる。

【図面の簡単な説明】

【図1】 本発明に係る一実施形態である連続音声認識装置のブロック図である。

【図2】 図1の連続音声認識装置における単語仮説絞込部6の処理を示すタイミングチャートである。

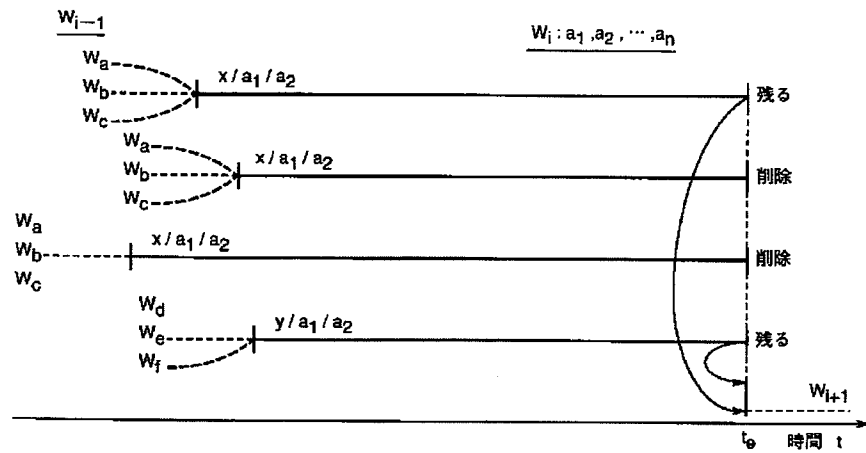
【図3】 図1の統計的言語モデル生成部20によって実行される統計的言語モデル生成処理を示すフローチャートである。

【図4】 図3のサブルーチンであるブラウンアルゴリズムを用いた分類処理を示すフローチャートである。

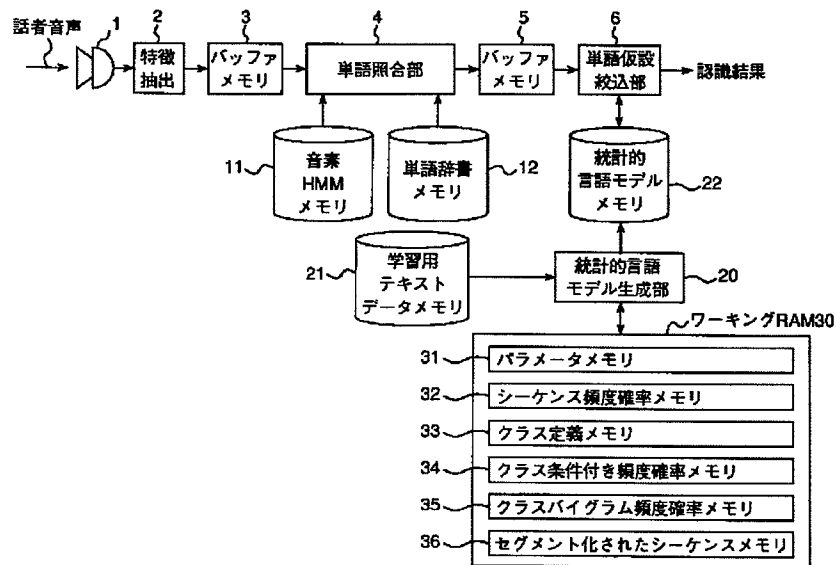
【符号の説明】

- 1…マイクロホン、
- 2…特徴抽出部、
- 3, 5…バッファメモリ、
- 4…単語照合部、
- 6…単語仮説絞込部、
- 11…音素HMMメモリ、
- 12…単語辞書メモリ、
- 20…統計的言語モデル生成部、
- 21…学習用テキストデータメモリ、
- 22…統計的言語モデルメモリ、
- 30…ワーキングRAM、
- 31…パラメータメモリ、
- 32…シーケンス頻度確率メモリ、
- 33…クラス定義メモリ、
- 34…クラス条件付き頻度確率メモリ、
- 35…クラスバイグラム頻度確率メモリ、
- 36…セグメント化されたシーケンスメモリ。

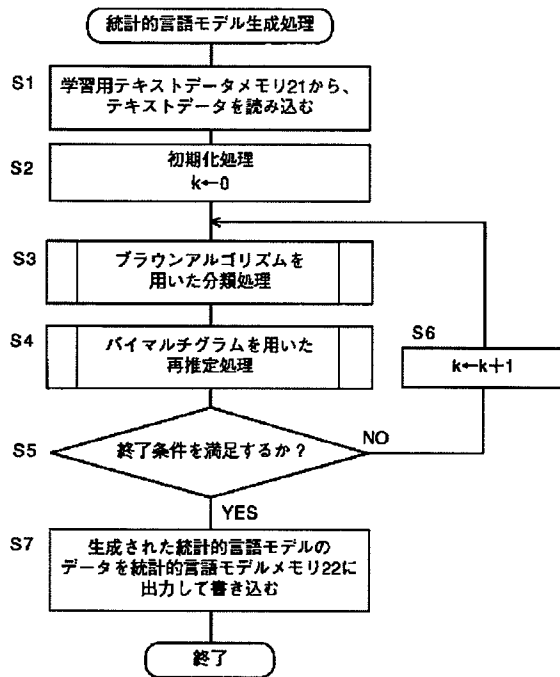
【図2】



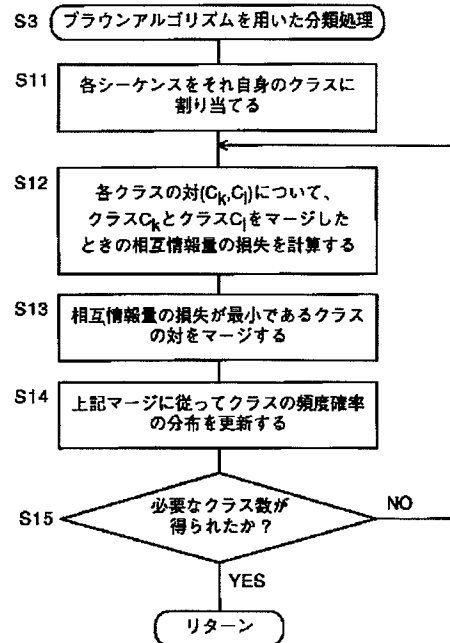
【図1】



【図3】



【図4】



フロントページの続き

(72)発明者 中嶋 秀治
京都府相楽郡精華町大字乾谷小字三平谷
5番地 株式会社エイ・ティ・アール音
声翻訳通信研究所内

(56)参考文献 DELIGNE S. "LANGUAGE MODELING BY VARIABLE LENGTH SEQUENCES: THEORETICAL FORMULATION AND EVALUATION OF MULTIGRAMS", ICASSP 1995, Vol. 1, pp169-172

Deligne S. "INFERENCE OF VARIABLE-LENGTH ACOUSTIC UNITS FOR CONTINUOUS SPEECH RECOGNITION", ICASSP 1997, Vol. 3, pp1731-1734

Frederic B. et. al. "Variable-Length Sequence Modeling: Multigrams", IEEE Signal Processing Letters, Vol. 2, No. 6, pp 111-113, JUNE 1995

(58)調査した分野(Int. Cl.⁷, DB名)

G10L 3/00 - 9/20

C12N 15/00

J I C S Tファイル (J O I S)